

# Anova (Analysis of Variance)

- M1:  $Y_i = \beta_0 + e_i$   $\xrightarrow{RSS}$   $\sum (Y_i - \bar{Y})^2 = SYY$   
 - M2:  $Y_i = \beta_0 + \beta_1 X_i + e_i$   $\xrightarrow{RSS}$   $SYY - \frac{SXY^2}{SXX}$

$$SS_{reg} = RSS_{M1} - RSS_{M2} = \frac{SXY^2}{SXX} = \begin{cases} \text{large} \rightarrow \text{favors M2} \\ \text{small} \rightarrow \text{favors M1} \end{cases}$$

$H_0$ : M1 is the right model  
 $H_A$ : M2 is the right model

$$SS_{reg} = \left( \frac{\sum (X_i - \bar{X}) Y_i}{\sqrt{SXX}} \right)^2 \sim \sigma^2 \chi_1^2$$

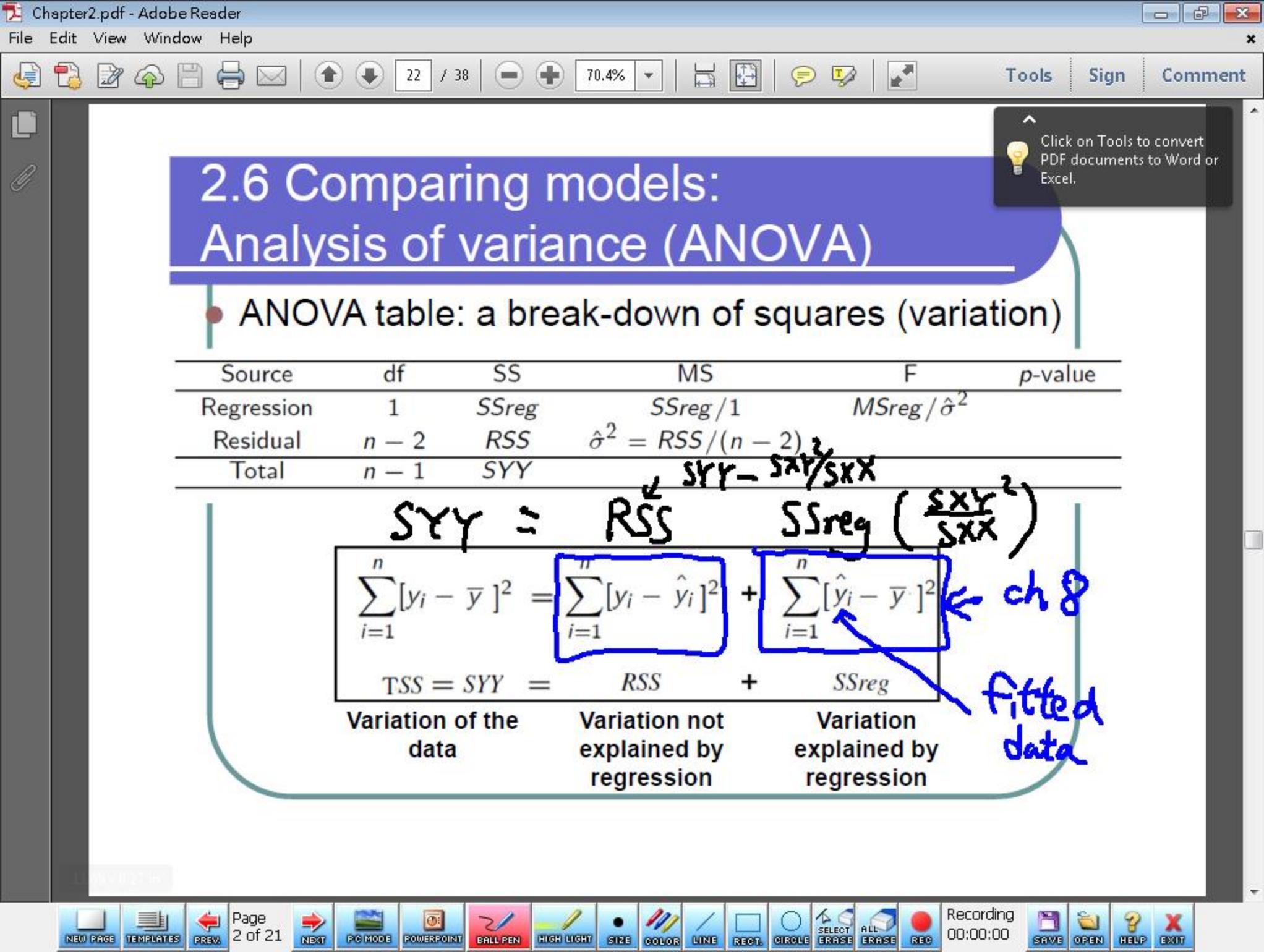
↑  
unknown

$$F\text{-stat} = \frac{SS_{reg}}{\frac{RSS}{n-2}} \sim \frac{\sigma^2 \chi_1^2}{\frac{\sum \hat{e}_i^2}{n-2} \sim \frac{\sigma^2 \chi_{n-2}^2}{n-2}} \sim F(1, n-2)$$

eg data  $\rightarrow$  F-stat = 3.4  $\xrightarrow{use R}$   
 If  $\alpha < 0.05 \Rightarrow$  reject  $H_0$   
 $\alpha \geq 0.05 \Rightarrow$  not reject  $H_0$



- Stat. facts
- ①  $\sum a_i Y_i \xrightarrow{CLT} \text{Normal}$
  - ②  $X_i \sim N(0, 1)$   
 $\Rightarrow X_1^2 + X_2^2 + \dots + X_k^2 \sim \chi_k^2$
  - ③  $X_1 \sim \chi_m^2$   
 $X_2 \sim \chi_n^2$   
 $\Rightarrow \frac{X_1/m}{X_2/n} \sim F(m, n)$
  - ④  $X_1 \sim N(0, 1)$   
 $X_2 \sim \chi_m^2$   
 $\Rightarrow X_1 / \sqrt{X_2/m} \sim t_m$



## 2.6 Comparing models: Analysis of variance (ANOVA)

- ANOVA table: a break-down of squares (variation)

Source	df	SS	MS	F	p-value
Regression	1	$SS_{reg}$	$SS_{reg}/1$	$MS_{reg}/\hat{\sigma}^2$	
Residual	$n - 2$	$RSS$	$\hat{\sigma}^2 = RSS/(n - 2)$		
Total	$n - 1$	$SYY$			

$$SYY = RSS + SS_{reg} \left( \frac{SXY^2}{SXX} \right)$$

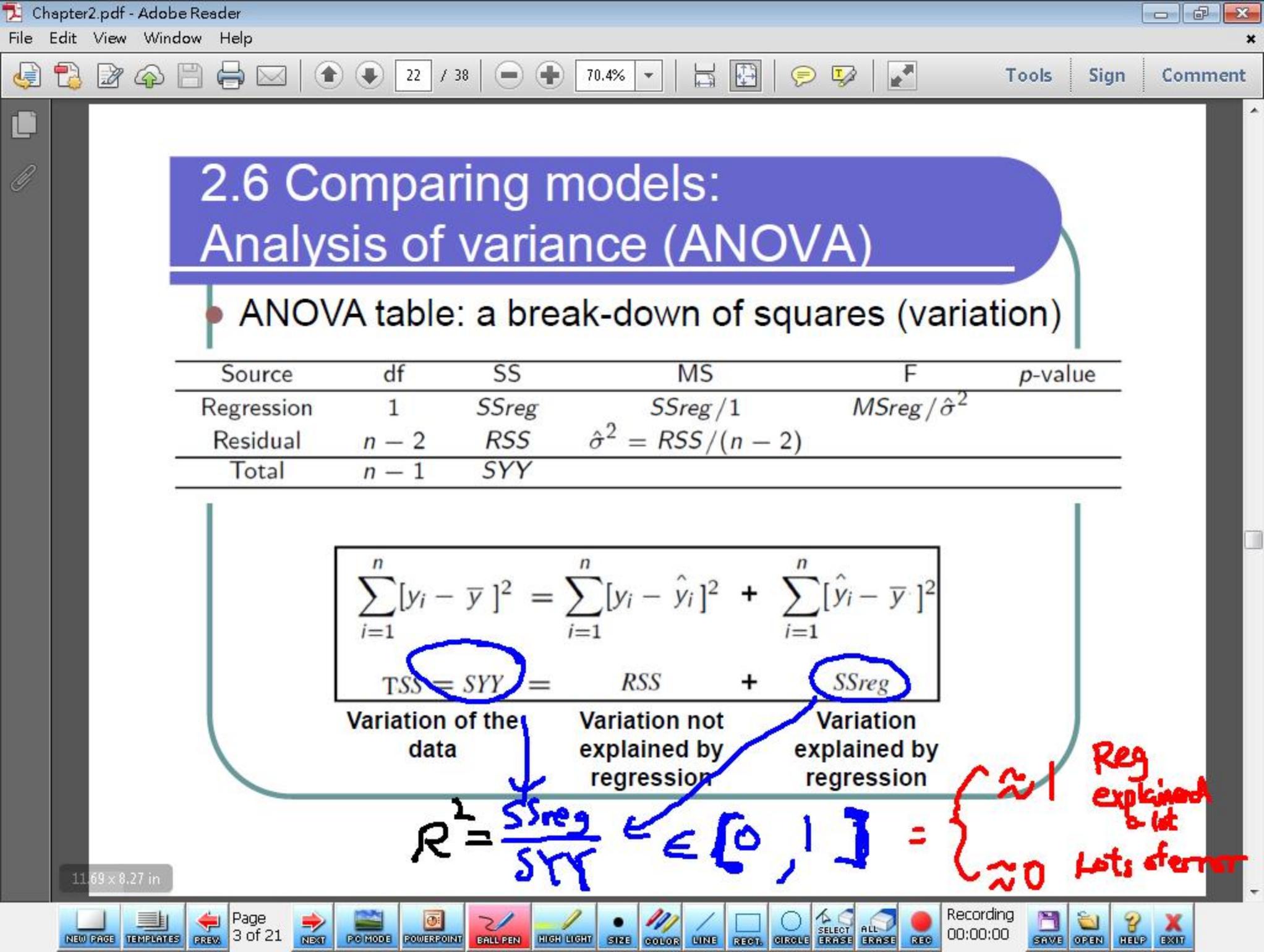
$$\sum_{i=1}^n [y_i - \bar{y}]^2 = \sum_{i=1}^n [y_i - \hat{y}_i]^2 + \sum_{i=1}^n [\hat{y}_i - \bar{y}]^2$$

TSS = SYY = RSS + SSreg

Variation of the data      Variation not explained by regression      Variation explained by regression

← ch 8  
 ← fitted data

Click on Tools to convert PDF documents to Word or Excel.



## 2.6 Comparing models: Analysis of variance (ANOVA)

- ANOVA table: a break-down of squares (variation)

Source	df	SS	MS	F	p-value
Regression	1	$SS_{reg}$	$SS_{reg}/1$	$MS_{reg}/\hat{\sigma}^2$	
Residual	$n - 2$	$RSS$	$\hat{\sigma}^2 = RSS/(n - 2)$		
Total	$n - 1$	$SYY$			

$$\sum_{i=1}^n [y_i - \bar{y}]^2 = \sum_{i=1}^n [y_i - \hat{y}_i]^2 + \sum_{i=1}^n [\hat{y}_i - \bar{y}]^2$$

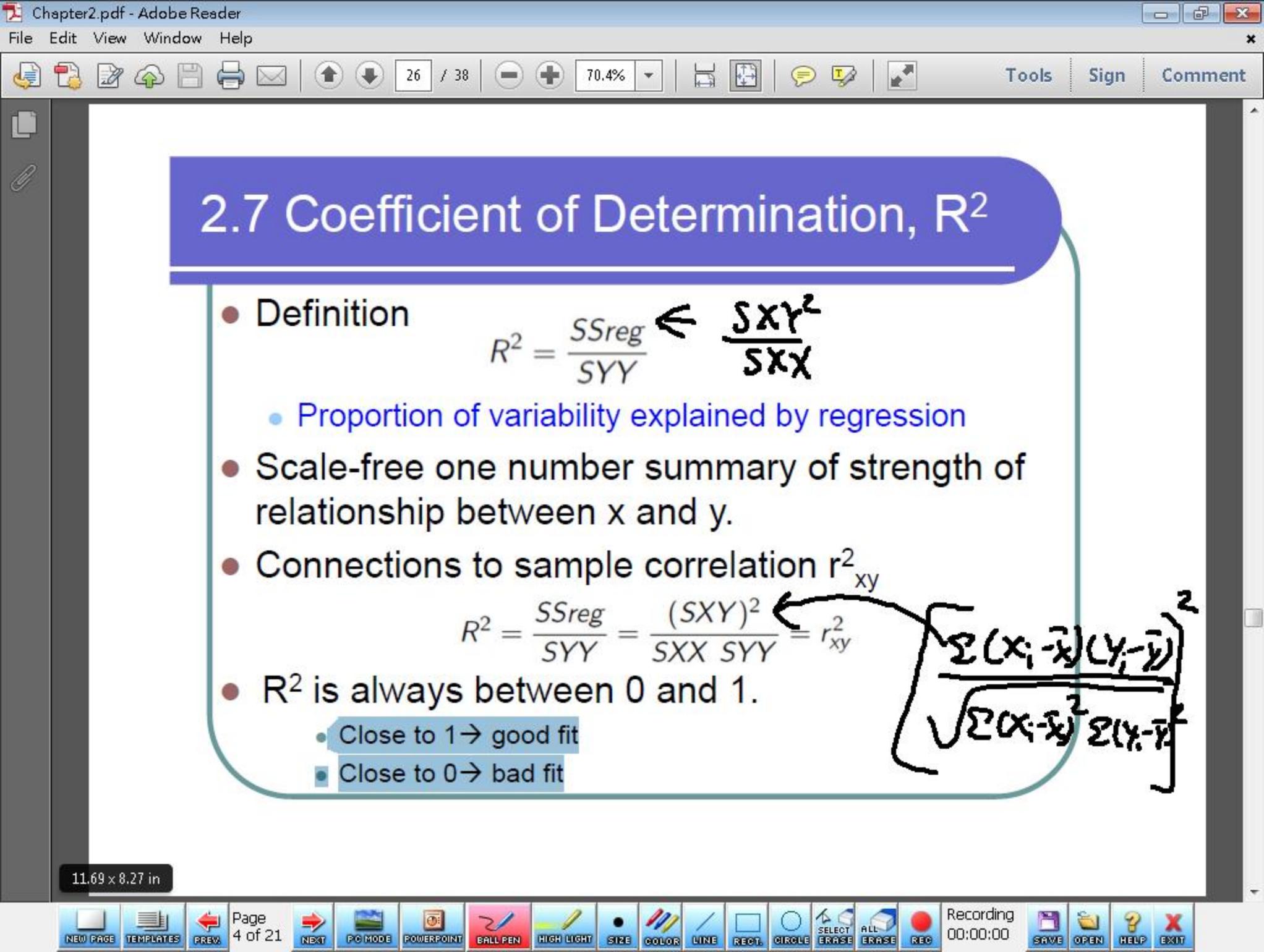
$TSS = SYY = RSS + SS_{reg}$

Variation of the data

Variation not explained by regression

Variation explained by regression

$R^2 = \frac{SS_{reg}}{SYY} \in [0, 1] = \begin{cases} \approx 1 & \text{Reg explained a lot} \\ \approx 0 & \text{Lots of error} \end{cases}$



## 2.7 Coefficient of Determination, $R^2$

- Definition

$$R^2 = \frac{SS_{reg}}{SYY} \leftarrow \frac{SXY^2}{SXX}$$

- Proportion of variability explained by regression
- Scale-free one number summary of strength of relationship between  $x$  and  $y$ .
- Connections to sample correlation  $r^2_{xy}$

$$R^2 = \frac{SS_{reg}}{SYY} = \frac{(SXY)^2}{SXX SYY} = r^2_{xy}$$

- $R^2$  is always between 0 and 1.
  - Close to 1  $\rightarrow$  good fit
  - Close to 0  $\rightarrow$  bad fit

$$\frac{\left[ \sum (x_i - \bar{x})(y_i - \bar{y}) \right]^2}{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}$$

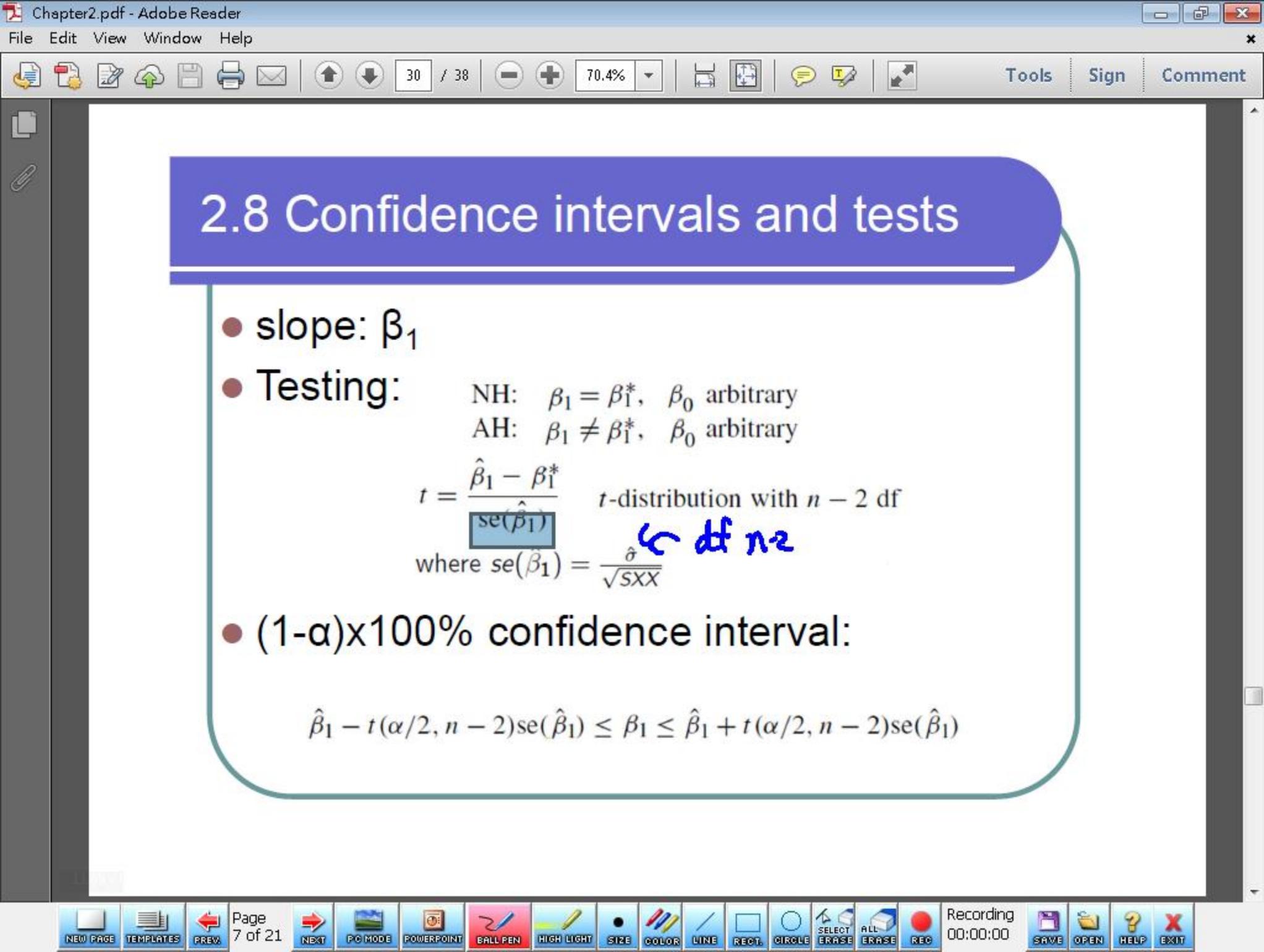
$$n=12$$

$$F = 2.4 = \frac{SS_{reg}}{RSS/10}$$

$$\begin{aligned} \text{Want } R^2 &= \frac{SS_{reg}}{S_{YY}} \\ &= \frac{SS_{reg}}{SS_{reg} + RSS} \\ &= \frac{SS_{reg}/RSS}{SS_{reg}/RSS + 1} \end{aligned}$$

$$S_{YY} = SS_{reg} + RSS$$

$$\therefore \frac{0.24}{0.24 + 1} = \frac{24}{124} = \frac{6}{31}$$



## 2.8 Confidence intervals and tests

- slope:  $\beta_1$
- Testing: NH:  $\beta_1 = \beta_1^*$ ,  $\beta_0$  arbitrary  
AH:  $\beta_1 \neq \beta_1^*$ ,  $\beta_0$  arbitrary

$$t = \frac{\hat{\beta}_1 - \beta_1^*}{\text{se}(\hat{\beta}_1)} \quad t\text{-distribution with } n - 2 \text{ df}$$

$$\text{where } \text{se}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{SXX}}$$

df n-2

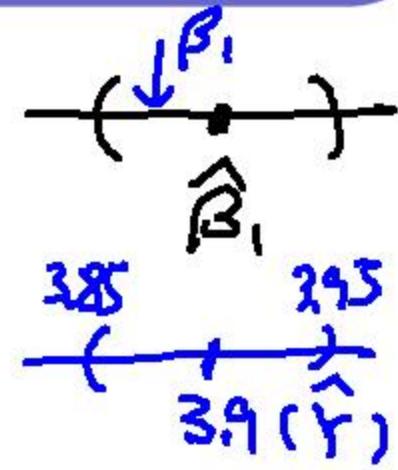
- $(1-\alpha) \times 100\%$  confidence interval:

$$\hat{\beta}_1 - t(\alpha/2, n - 2)\text{se}(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + t(\alpha/2, n - 2)\text{se}(\hat{\beta}_1)$$

$Y_i = \beta_0 + \beta_1 X_i + e_i$  ← individual variation

## 2.8 Confidence intervals and tests

- Regression model:
  - $E(Y|X=x) = \beta_0 + \beta_1 x$
- Quantities of interests



- Intercept:  $\beta_0$
- Slope:  $\beta_1$
- Prediction: If we observe  $x_*$ , what is the  $y$ ?
- Fitted value:  $E(Y|X=x)$  for different values of  $x$

individual  
average

- Confidence intervals give estimates for the above quantities of interests

Remember STAT2001  $(X_1, \dots, X_n) \sim$  true mean  $\mu$

$$H_0: \mu = 3$$

$$H_A: \mu \neq 3$$

Testing

C.I. for  $\mu$

$$t = \frac{\bar{x} - 3}{\sqrt{\text{Var}(\bar{x} - 3)}} = \frac{\bar{x} - 3}{\hat{\sigma}/\sqrt{n}}$$

$$\sim \frac{\text{Normal}}{\sqrt{\chi^2}} \sim t\text{-dist}$$



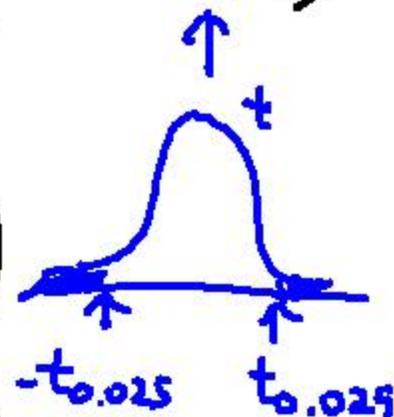
$$0.05 = P\left(-t_{0.025} < \frac{\bar{x} - \mu}{\hat{\sigma}/\sqrt{n}} < t_{0.025}\right)$$

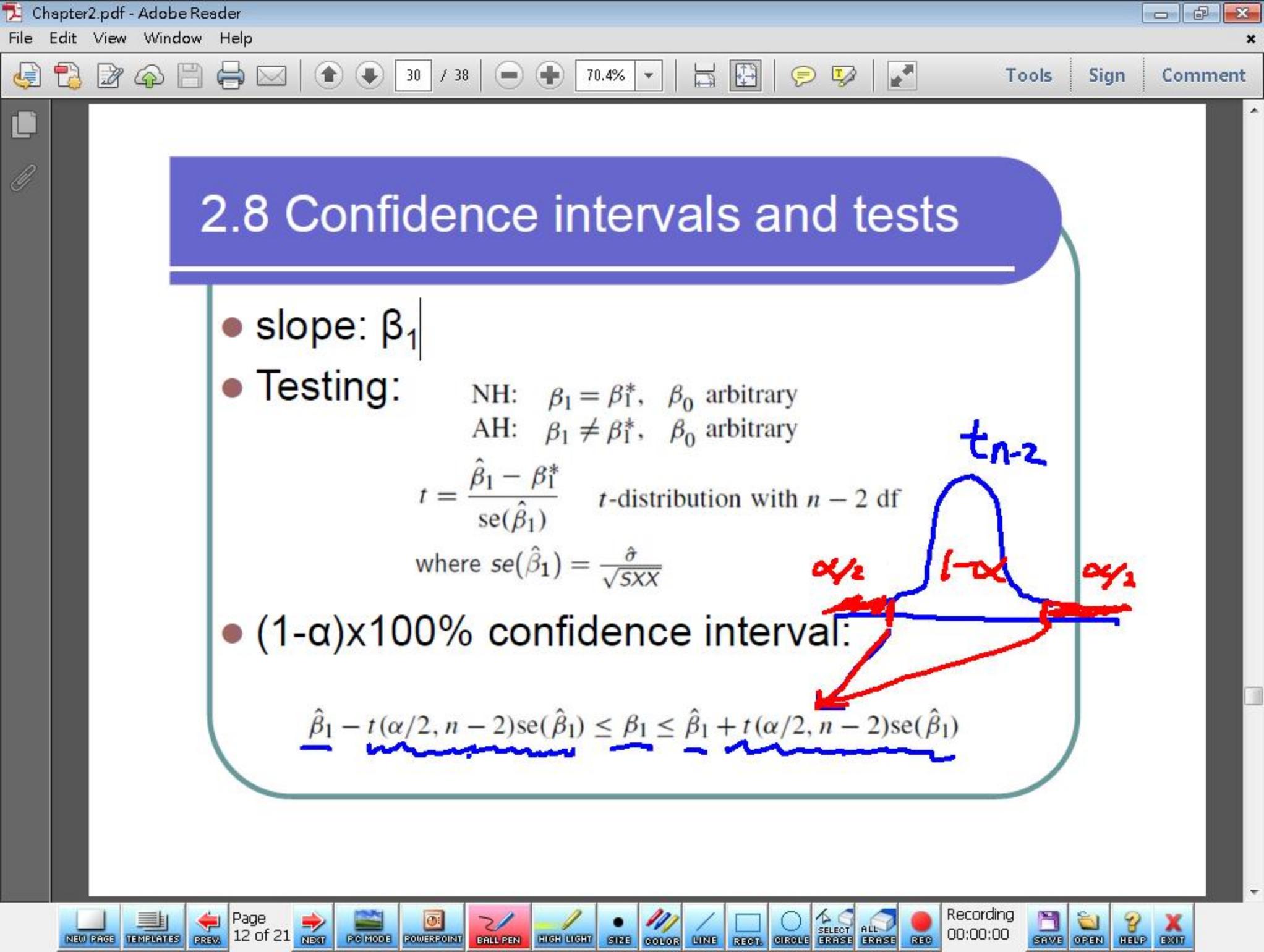
$$\Rightarrow 0.05$$

$$= P\left(\bar{x} - t \frac{\hat{\sigma}}{\sqrt{n}} < \mu < \bar{x} + t \frac{\hat{\sigma}}{\sqrt{n}}\right)$$

$$\Rightarrow \mu \in \left(\bar{x} - t \frac{\hat{\sigma}}{\sqrt{n}}, \bar{x} + t \frac{\hat{\sigma}}{\sqrt{n}}\right)$$

C.I.





## 2.8 Confidence intervals and tests

- slope:  $\beta_1$
- Testing: NH:  $\beta_1 = \beta_1^*$ ,  $\beta_0$  arbitrary  
AH:  $\beta_1 \neq \beta_1^*$ ,  $\beta_0$  arbitrary

$$t = \frac{\hat{\beta}_1 - \beta_1^*}{\text{se}(\hat{\beta}_1)} \quad t\text{-distribution with } n - 2 \text{ df}$$

$$\text{where } \text{se}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{SXX}}$$

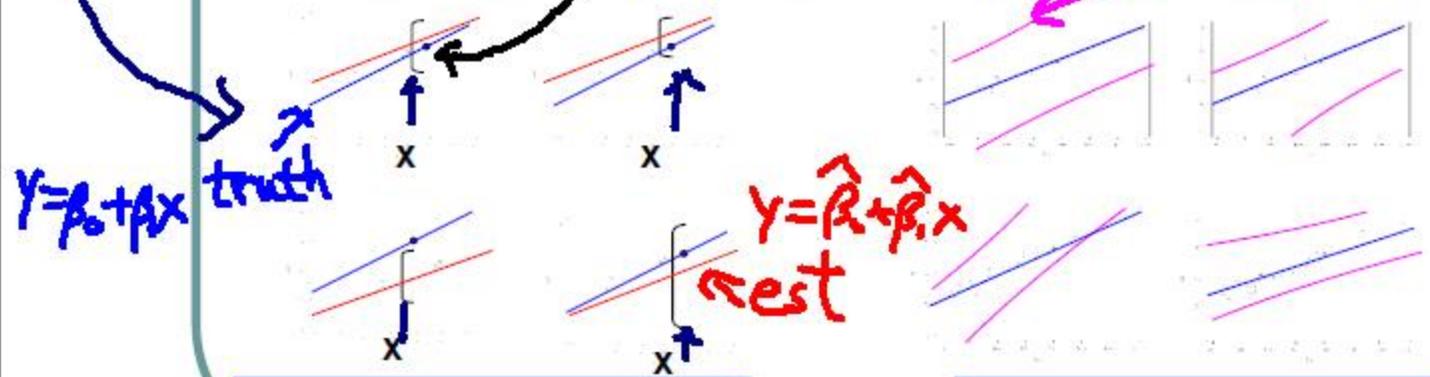
- $(1-\alpha) \times 100\%$  confidence interval:

$$\hat{\beta}_1 - t(\alpha/2, n - 2)\text{se}(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + t(\alpha/2, n - 2)\text{se}(\hat{\beta}_1)$$

$$\begin{aligned}
& \text{Var} \left( (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) X_* + e \right) \\
= & \text{Var}(\beta_0 - \hat{\beta}_0) + X_*^2 \text{Var}(\beta_1 - \hat{\beta}_1) + \underbrace{\sigma^2}_{\text{individual}} \\
& + 2 \text{Cov}(\beta_0 - \hat{\beta}_0, \beta_1 - \hat{\beta}_1) \\
& + \cancel{2 \text{Cov}(\beta_0 - \hat{\beta}_0, e)} \quad \begin{array}{l} \text{from new obs} \\ \text{from past data} \end{array} \\
& + \cancel{2 \text{Cov}(\beta_1 - \hat{\beta}_1, e)}
\end{aligned}$$

# 2.8 Confidence intervals and bands

- Confidence interval (at each point  $x$ )
  - For each of  $x$ ,  $P(E(Y|X=x) \text{ in C.I.}) = 1-\alpha$
- Confidence band (for the entire line)
  - $P(\text{For all } x, E(Y|X=x) \text{ in C.B.}) = 1-\alpha$  **C.B.**



For  $n$  C.I.s,  $n(1-\alpha)$  of them covers the true value at  $x$

For  $n$  C.B.s,  $n(1-\alpha)$  of them covers the whole true regression line